



2023 OFA Virtual Workshop

MULTICAST CONGESTION CONTROL SUPPORT IN ROCE V2

Christoph Lameter, Senior IT Experte

Deutsche Börse AG



Multicast Congestion Control support in RoceV2

christoph@lameter.com

March 22nd, 2023

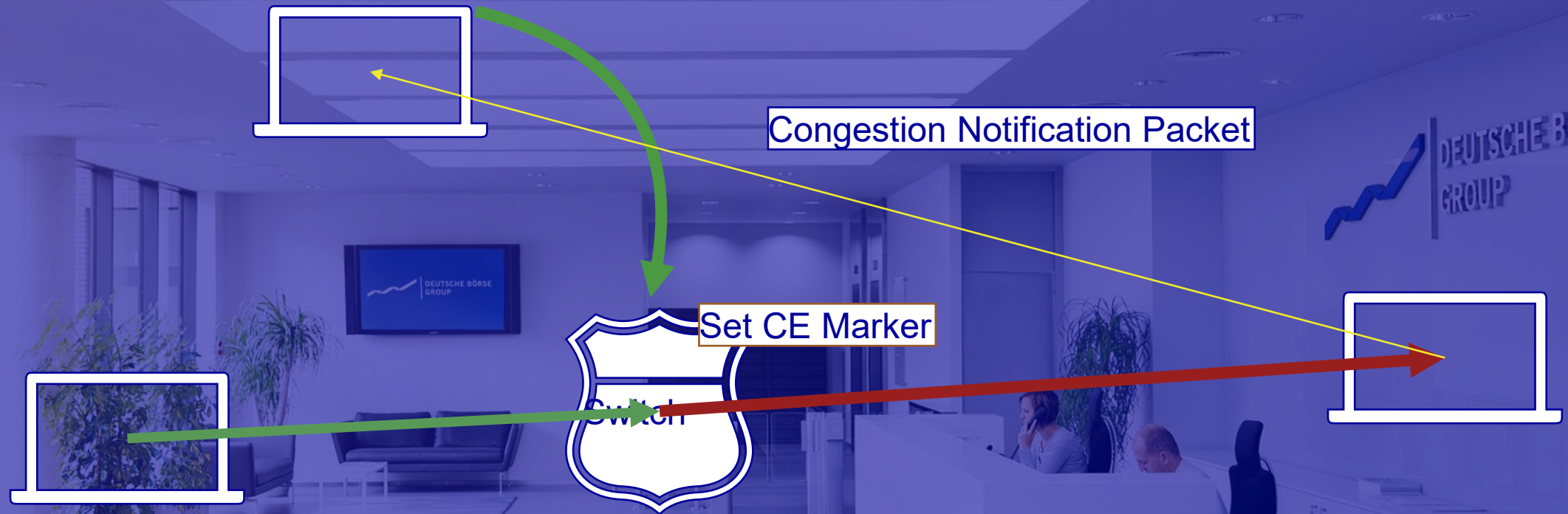


DEUTSCHE BÖRSE
GROUP

Index

- Page 3 -> Roce V2 Unicast congestion Control
- Page 5 -> Multicast: Magic in Networking
- Page 7 -> Current Handling of Multicast Congestion
- Page 9 -> Proposal on how to manage Multicast Congestion
- Page 11 -> IETF Standards and marking Multicast Traffic with the CE bits
- Page 13 -> CNP unicast as a result of ROCE multicast packets with CE set
- Page 15 -> CNP Flooding Dangers due to sending multicast with CE bit set
- Page 17 -> Flow Control for Multicast (IEEE 802.1Qbb)
- Page 19 -> Test status RoceV2 congestion control
- Page 22 -> Testing Tools (ib2roce, mcsend, mclisten)
- Page 24 -> Deutsche Boerse AG Legalese and Corporate Slides

Roce V2 Unicast Congestion Control



IBTA Spec A17: Roce V2 Congestion Control



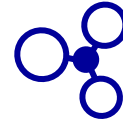
ECN / CE Signalling

A17.3.1.1.3 / 17.9.3 specifies the use of the ECN field to indicate congestion in the switch to the receiver. The receiver can then send CNP packets to slow down the sender. Support for ECN / CE is optional (CA17-5).



CNP packets slow down senders

A receiver sends CNPs to the individual sender if the CE flag is set in the header (CA17-45, CA17-44). Switches have advanced strategies to mark packet more or less frequently depending on the level of congestion to trigger CNPs.



Link-Layer Flow-Control IEEE802.1Qbb

Roce may also use Link Layer flow control (Pause Frames, PFC) in addition to CNP packets. This is a defined flow control method at the *Ethernet Layer* that ROCE packets can make use of. The sending rate on the one side of a cable is slowed down to mitigate network congestion.

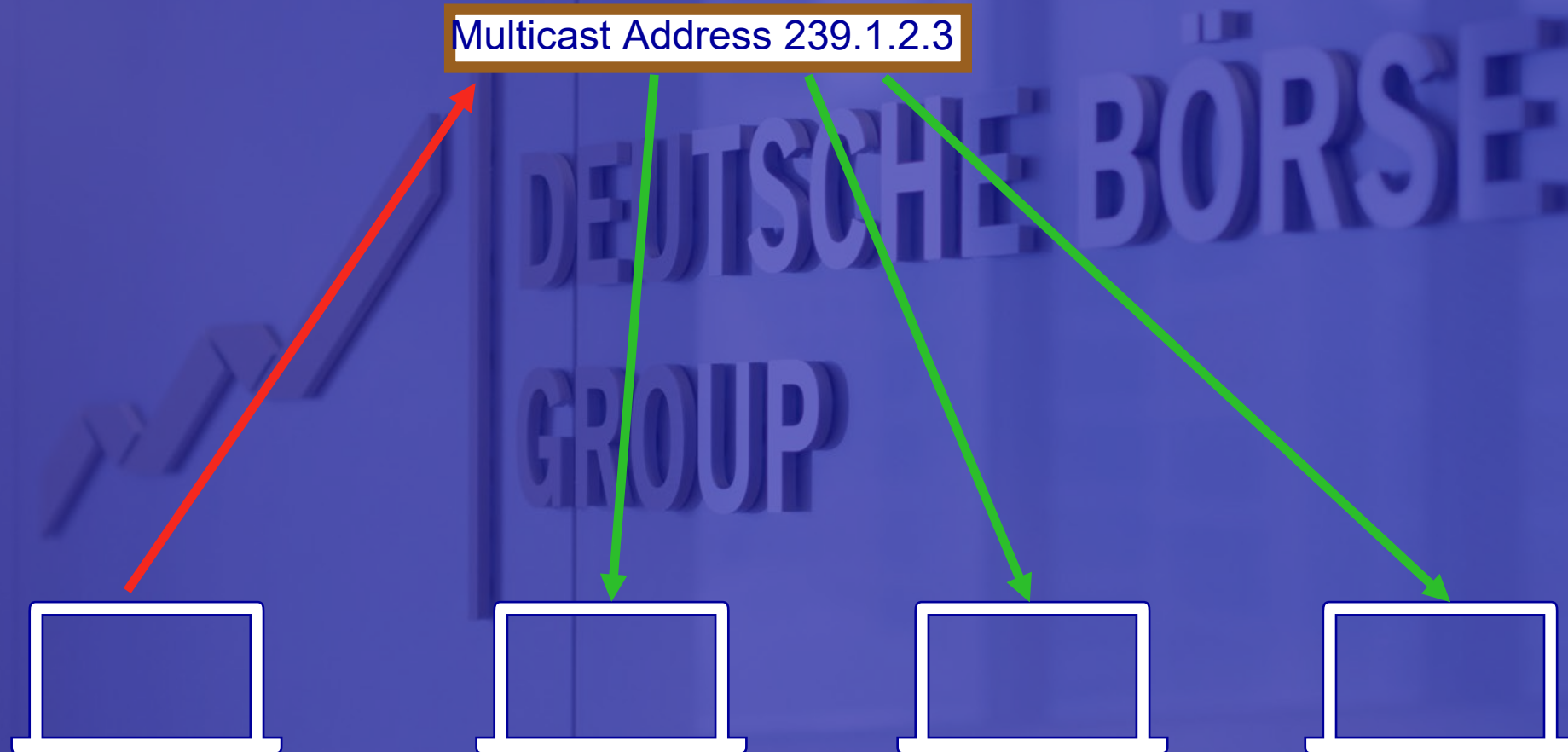


Lossless operations A17.9.1.

ROCE handling should be lossless. This is in practice realized in the following way:

1. CNP packets to slow down the sender
2. Activation of Link Layer Flow control if CNPs are not sufficient.

Multicast : Magic in Networking



Multicast: IP addresses that one can subscribe to and that can reach multiple recipients with one datagram

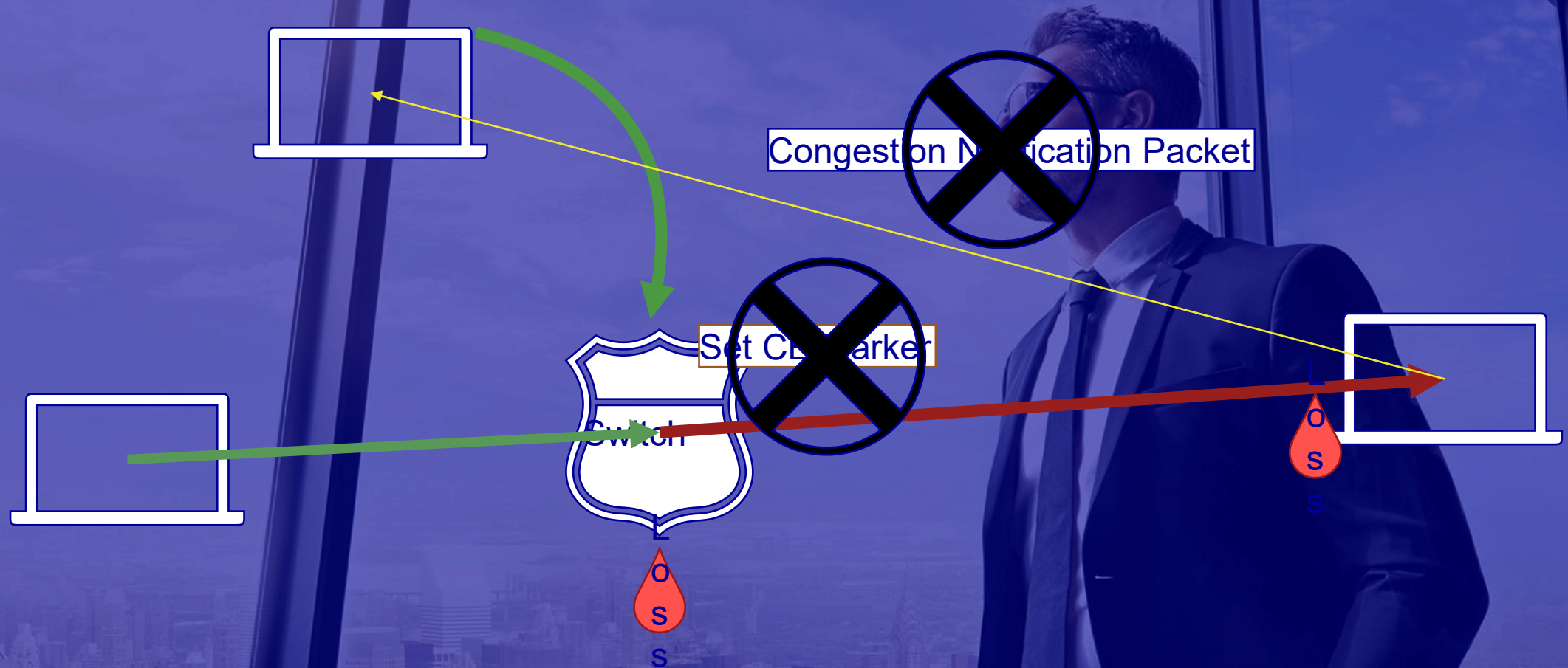
Multicast is a widely used standard to implement a publish/subscribe mechanism which has been available since the 1970s and was widely adopted in particular in the financial services industry in the 1980s (see RFC 1112). It is supported on Ethernet and IP based networks. Multicast is also available on Infiniband (IB Specification Volume 1, 3.5.11) and existing implementations have RDMA APIs to manage Multicast subscriptions (*rdma_join_multicast()*, *rdma_leave_multicast()*).

A special *IPv4 address range* 224.0.0.0 – 239.255.255.255 is reserved for multicast. Subscribers can join by specifying the multicast address and will then receive messages sent with the destination of that multicast address. There are a variety of popular protocols to optimize the management and routing of traffic to multicast recipients like **IGMP** (RFC 3376) , **PIM Sparse/Dense** (RFC 2362). In Infiniband the subnet manager is responsible to optimize multicast traffic.

The IP multicast range has an associated MAC address range when used in Ethernet (IPv4 = 01-00-5e-xx-xx-xx). A translation of Multicast IP to MAC addresses is easily possible via a simple calculation and preserves the lower bits of the IP address in the MAC address.

In IPv6 a prefix is reserved for multicast (0xFF00::/8). Infiniband adopts the IPv6 conventions for the specification of multicast groups in the form of a special GID the **MGID**. MGIDs are associated with **MLIDs** by the Subnetmanager for routing in the Infiniband switches. Infiniband switches are then able to do cut-through of a single incoming datagram message to multiple output ports. This process is typically resulting in packet delivery in under a microsecond to a large number of end points. The rapid event publication via this mechanism is one of the reasons for the use of RDMA in the financial services industry.

Current Handling of Multicast Congestion

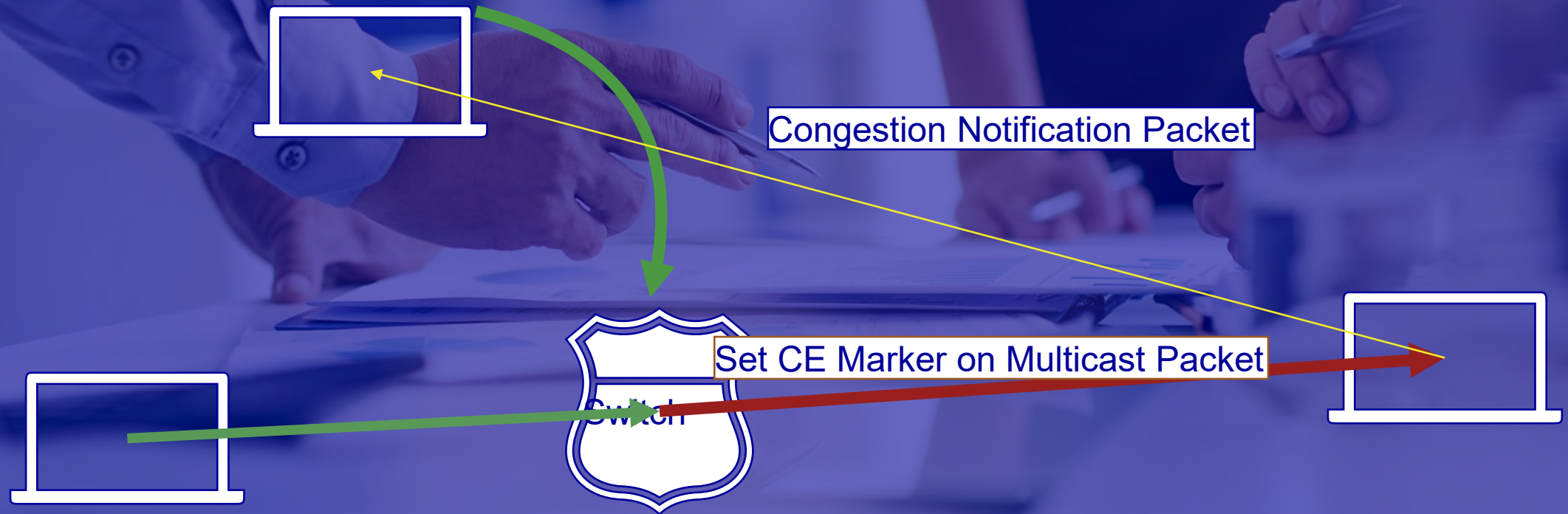


Current handling of Roce V2 multicast traffic

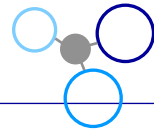
- Link Level flow control is possible if the switch supports in (one known vendor). PFC frames will slow down all traffic on the link. Otherwise traffic will be dropped.
- CE bit is not set by the switches if links become congested since switches do not support congestion control for Multicast.
- Therefore the NICs do not send CNPs.
- Traffic on the link does not slow down.
- Massive packet loss is possible which leads to unrecoverable failure of applications expecting reliable multicast solutions as provided by middleware vendors.
- Multicast based RDMA applications work reliably on Infiniband but not on ROCE.
- This is contrary of A17.9.1 which states that ROCE v2 handling should be lossless. However, that statement seems to be only be true for unicast traffic. Annex 17 does not make a distinction between multicast and unicast and does not discuss the issues related to multicast congestion control.
- Apart from one vendors support for Multicast Flow Control, currently available switches have no option of making ROCE multicast reliable. The one vendor does not support CNPs so all traffic coming from an interface has to be throttled instead of only the stream that causes the overload.



Proposal on how to manage Multicast Congestion



Multicast Congestion Control / Testing of Multicast?



IBTA Spec proposal

A17.3.1.1.3 ECN

- Add the following paragraph:
- **Note that the use of ECN/Congestion management is not envisioned to be restricted to unicast packets. RFC 3168 is considering the use of ECN also for multicast. Multicast packet congestion can be treated in the same way as unicast traffic. For additional considerations see Section 17.9.3.**

A17.9.3

- Add the following paragraph:
- **Note that congestion control through ECN or Flow-Control may be implemented for multicast packets so that reliable multicast is possible like under Infiniband. For the purpose of congestion control the multicast address is treated as a single destination like a regular unicast address. Therefore, congestion on the path to a single receiver of the multicast group may slow down the reception of traffic for all receivers (and then be compatible with the way Infiniband handles congestion).**
- **Switches and NICs typically allow a fine grained QoS configuration which allows customization of congestion handling. This customization is possible using the existing QoS configuration API that is familiar to network administrators if ECN handling and Ethernet Flow Control is not restricted to unicast packets by the switches and NICs.**

IETF Standards and marking multicast traffic with the CE bits



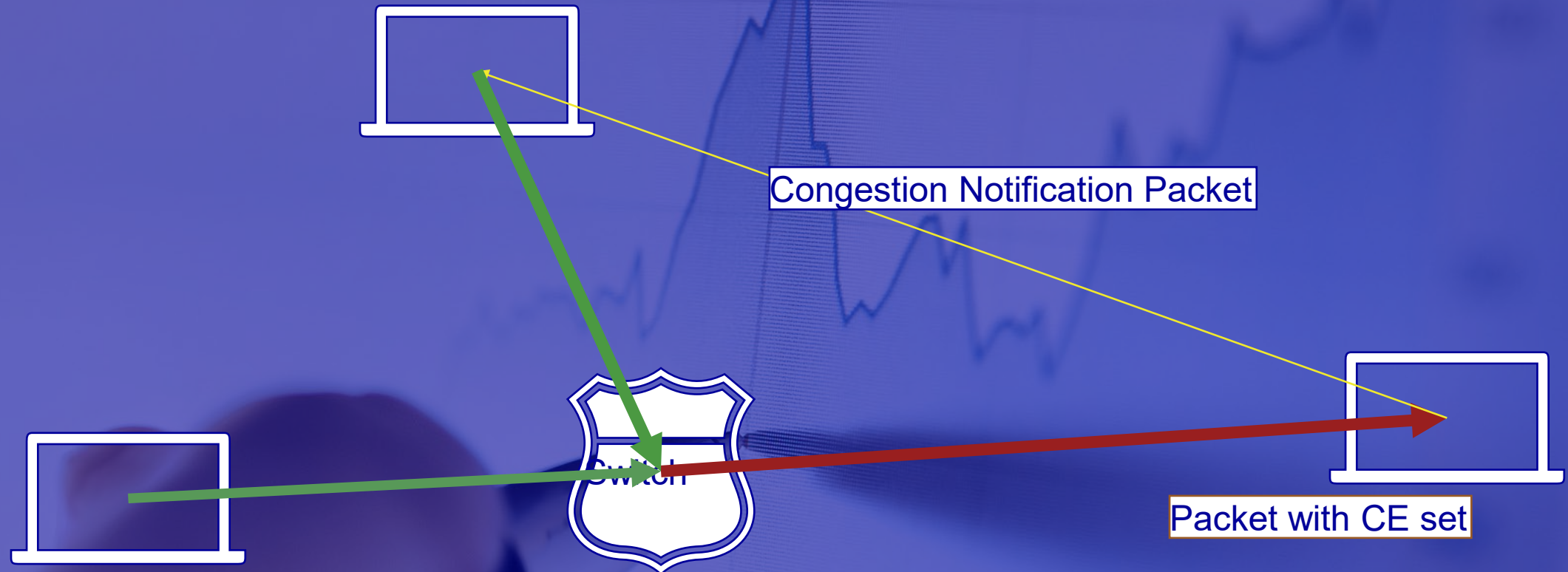
According to A17.9.3 Congestion management relies on “*the mechanism defined in*” RFC3168 (ECN). RFC 3168 is primarily concerned with defining a congestion control mechanism for the TCP protocol by allowing detection of congestion before queues overflow. ECN is relying on the transport protocol (Section 6) to allow a reaction to packets that have encountered congestion. The receiver notices that the CE bits are set and performs a transport specific reaction to reduce congestion in the fabric.

RFC 3168 envisions other transport protocols to be using this mechanism mentioning among others “*unreliable and reliable multicast transport protocols*” (Section 6).

ROCE V2 is such a transport protocol where the IBTA has defined congestion control mechanism through the congestion notifications packets (CNP). These have been so far be seen as restricted to unicast. It looks like it is within the intended scope of RFC 3168 if the use of CNP would be extended to cover multicast as well.

The IBTA spec does not restrict the use of CNPs to unicast. However, the actual implementations by the hardware vendors have so far restricted the use of CNPs to unicast.

CNP unicast as a result of ROCE multicast packets with CE set



CNP replies as a result of receiving Multicast Packets marked with CE

A17.9.3 states:

CA17-44: If RoCEv2 Congestion Management is supported, upon receiving a valid RoCEv2 packet with a value of '11 in its IP.ECN field the HCA shall generate a RoCEv2 CNP formatted as shown in Figure 360 on page 2000 directed to the source of the received packet. The HCA may choose to send a single CNP for multiple such ECN marked packets on a given QP.

So far, the CNPs are only sent as a result of receiving unicast ROCE v2 packets. The use of this mechanism when receiving a multicast packet is straightforward. A unicast CNP packet is sent when a multicast packet with the CE bits set is detected. The destination is the sender of the multicast traffic which is the source address of the multicast packet (which is a unicast address).

The sender will then receive the unicast CNP packet and moderate the output of the multicast QP to reduce congestion on the fabric.

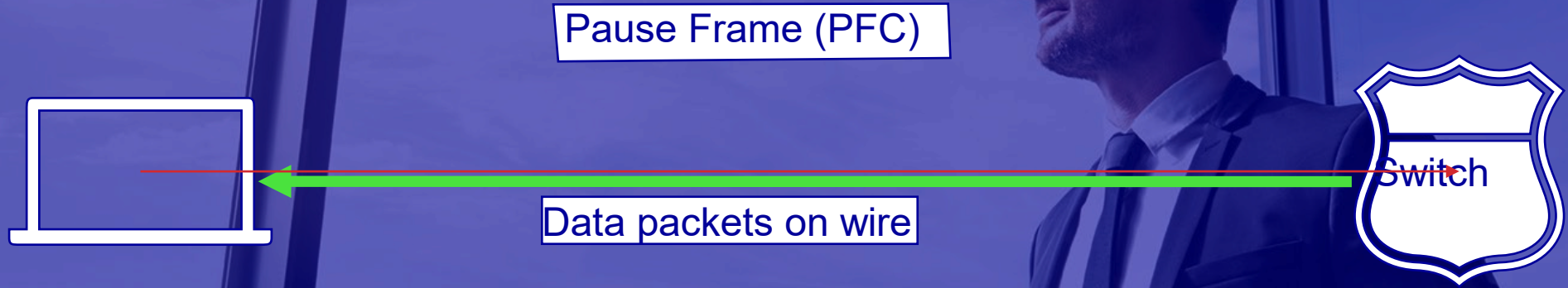
CNP Flooding Dangers due to sending multicast with CE bit set



Multicast replication can result in single packets marked with the CE bit arriving at multiple endpoints which may all send CNP responses to the sender. It is therefore advisable to be able to limit the number of CNPs in a certain time interval. This has already been an issue for unicast and therefore these throttling mechanisms already exist at the sender/receiver as well as the switches in the hardware known to the author. These are currently not used for multicast. There are 3 points the network where these flooding issues could be handled:

1. The receiver may receive too frequent CE packets and respond to reach with CNP. These issues have already been observed with unicast and receivers limit the number of responses to CE packets already. Switches have the ability to control the probabilities of packets being marked with CE already.
2. Switches may receive too many CNPs. The switch can reduce the number of packets marked as CE to reduce CNPs to a desirable rate. The rate and logic of marking multicast packets as CE may differ from unicast packets and that may require switch vendors to introduce new tuneable parameters. Switches may use their knowledge of multicast replication trees to provide advanced means of congestion control.
3. Sender may receive too many CNPs. The sender may limit the number of CNPs that it reacts to in a given timeframe. It seems that there are already heuristics in place to avoid excessive throttling of senders for unicast.

Flow Control for Multicast



IEEE 802.1Qbb considers any traffic between the two endpoints of a link. The standard is therefore not concerned with unicast versus multicast. The PFC frame itself is a multicast frame through. However, that special frame is only used for communication between the two endpoints of the link. The PFC frame is never forwarded to other ports of the switch. The standard establishes priority classes for PFC support (defined in IEEE 802.1Q) and allows the slowing down of senders for packets of a particular category.

This means that if unicast and multicast packets are treated the same then the mechanism works for either type of packet.

This is the case for switches from Nvidia. However, Cisco switches do not apply priorities to multicast traffic. Cisco has build a special queue for multicast traffic in which all multicast traffic accumulates bypassing the queues for the traffic classes. There are no provisions for congestion support for this special queue. This then defeats the mechanism envisioned in IEEE 802.1Qbb. Cisco switches not send PFC frames if multicast traffic is congested.

RDMA NICs maybe configured send PFCs to the switch if the applications are slow in picking up the traffic. This does work with the Cisco switches resulting in traffic to accumulate on the switches. The Cisco switch will have to drop packets when too many packets accumulate since it is not able to send PFCs to the sender of the multicast traffic.

Test status RoceV2 congestion control

Known switch Support for Congestion Management in ROCE V2

Feature	Unicast PFC	Unicast CNP	Multicast PFC	Multicast CNP
Cisco	Yes	Yes	Does not allow classification of multicast traffic into IEEE802.1Q priority classes.	Tests fail. ^[SEP] “Some ASICs cant do it.”
Nvidia	Yes	Yes	Yes	Tests fail. Working on getting this implemented.
Arista	TBD	TBD	“Should work” (David Snowdon)	TBD

Known NIC Support for ROCE v2 Multicast

Feature	Receive	Loopback Send	Remote Send
Nvidia	Yes	Yes	Yes
Intel	Yes	Yes	Patch for testing against upstream drivers exist to enable multicast send
Broadcom	No MC support in Linux source tree	-	-

Test tools and methods

We experienced some issues with the standard Infiniband testing tools (**ib_send_XX**) typically used for testing RDMA packet flow:

1. Multicast join and leave operations are not implemented in many tools.
2. In ROCE v2 mode the testing tools fall back to *broadcasting* on an IPv6 address instead of performing multicast.
3. Multicast through rdma_cm libraries (as used by middleware applications) is *not supported* by most tools

In order to have some tools that work with the RDMACM libraries we wrote some tools called **mclisten** and **mcsender**. These are simple to operate and provide a way to test in the same way that our middleware operates with the RDMA APIs.

These tools are available as a upcoming contribution to the rdma-core package on github. See <https://github.com/clameter/rdma-core>. Check out the branch ib2roce.

Corporate Slides

As an international exchange organisation and innovative market infrastructure provider, Deutsche Börse Group offers its customers a wide range of products, services and technologies covering the entire value chain of financial markets. It organises markets characterised by integrity, transparency and safety for investors who invest capital and for companies that raise capital.

Its business areas include the provision of index and ESG data, analytics and research solutions, trading and clearing services for investment instruments, securities settlement and custody, collateral and liquidity management, and investment fund services. In addition, the Group develops state-of-the-art IT solutions and offers IT systems all over the world.

With over 10,500 employees, the company has its headquarters in the financial centre of Frankfurt/Main, as well as a strong global presence in Luxembourg, Prague, London and Zug, in New York and Chicago, in Hong Kong, Singapore, Beijing, Tokyo and Sydney – and at other locations for its customers all over the world.

Contact

Corporate & Corporate Engagement

Deutsche Börse AG
Mergenthalerallee 61
65760 Eschborn

E-mail corporate.communications@deutsche-boerse.com

Disclaimer

© Deutsche Börse Group 2022

This publication is for informational purposes only. None of the information in this publication constitutes investment advice and does not constitute an offer to sell or a solicitation of an offer to purchase any contract, share or other financial instrument. This publication is not intended for solicitation purposes but only for use as general information. All descriptions, examples and calculations contained in this publication are for illustrative purposes only.

Deutsche Börse AG, Frankfurter Wertpapierbörse (FWB®, the Frankfurt Stock Exchange), Eurex Frankfurt AG, Eurex Deutschland and Eurex Clearing AG do not represent that the information in this publication is comprehensive, complete or accurate and exclude liability for any consequence resulting from acting upon the contents of this or another webpublication, in so far as no wilful violation of obligations took place or, as the case may be, no injury to life, health or body arises or claims resulting from the Product Liability Act are affected.

Securities traded on the Frankfurt Stock Exchange and Eurex derivatives (other than EURO STOXX 50® Index Futures contracts, EURO STOXX® Select Dividend 30 Index Futures contracts, STOXX® Europe 50 Index Futures contracts, STOXX® Europe 600 Index Futures contracts, STOXX® Europe Large/Mid/Small 200 Index Futures contracts, EURO STOXX® Banks Sector Futures contracts, STOXX® Europe 600 Banks/Industrial Goods & Services/Insurance/Media/Personal & Household Goods/Travel & Leisure/Utilities Futures contracts, Dow Jones Global Titans 50 IndexSM Futures contracts, DAX® Futures contracts, MDAX® Futures contracts, TecDAX® Futures contracts, SMIM® Futures contracts, SLI Swiss Leader Index® Futures contracts, RDXxt® USD - RDX Extended Index Futures contracts, Eurex inflation/commodity/weather/property and interest rate derivatives) are currently not available for offer, sale or trading in the United States nor may they be offered, sold or traded by persons to whom US tax laws apply.

The fund shares listed in XTF Exchange Traded Funds® are admitted for trading on the Frankfurt Stock Exchange. Users of this information service who legally reside outside Germany are herewith advised that sale of the fund shares listed in XTF Exchange Traded Funds may not be permitted in their country of residence. The user makes use of the information at their own risk.

Legal validity of this disclaimer

In the event that individual parts of or formulations contained in this text are not, or are no longer, legally valid (either in whole or in part), the content and validity of the remaining parts of the document are not affected.

Trademarks

The following names and designations are registered trademarks of Deutsche Börse AG or an affiliate of Deutsche Börse Group:

1585®; A7®; Buxl®; C7®; CDAX®; CEF®; CEF alpha®; CEF ultra®; CFF®; Classic All Share®; Clearstream®; CX®; D7®; DAX®; DAXglobal®; DAXplus®; DB1 Ventures®; DBIX Deutsche Börse India Index®, Deutsche Börse®; Deutsche Börse Capital Markets Partner®; Deutsche Börse Commodities®; Deutsche Börse Venture Network®; Deutsches Eigenkapitalforum®; DivDAX®; eb.rexx®; eb.rexX Jumbo Pfandbriefe®; ERS®; eTriParty®; Eurex®; Eurex Bonds®; Eurex Clearing Prisma®; Eurex Improve®; Eurex Repo®; Euro GC®; ExServes®; EXTF®; F7®; FDAX®; FWB®; GC Pooling®; GCPI®; GEX®; Global Emission Markets Access – GEMA®; HDAX®; iNAV®; L-DAX®; L-MDAX®; L-SDAX®; L-TecDAX®; M7®; MDAX®; N7®; ODAX®; ÖkoDAX®;PROPRIS®; REX®; RX REIT Index®; Scale®; SCHATZ-FUTURE®; SDAX®; ShortDAX®; StatistiX®; T7®; TecDAX®; Technology All Share®; TRICE®; USD GC Pooling®; VDAX®; VDAX-NEW®; Vestima®; Xscreen®, Xemac®; Xentric®, Xetra®; Xetra-Gold®; Xpect®; Xpider®, XTF®; XTF Exchange Traded Funds®; We make markets work®

The names and trademarks listed above do not represent a complete list and, as well as all other trademarks and protected rights mentioned in this publication, are subject unreservedly to the applicable trademark law in each case and are not permitted to be used without the express permission of the registered owner. The simple fact that this publication mentions them does not imply that trademarks are not protected by the rights of third parties.

The STOXX® indices, the data included therein and the trademarks used in the index names are the intellectual property of STOXX Ltd., Zug, Switzerland and/or its licensors. Eurex' derivatives based on the STOXX indices are in no way sponsored, endorsed, sold or promoted by STOXX and its licensors and neither STOXX nor its licensors shall have any liability with respect thereto.

STOXX iSTUDIO® is a registered trademark of STOXX Ltd., Zug, Switzerland.

TRADEGATE® is a registered trademark of Tradegate AG Wertpapierhandelsbank.

EEX® is a registered trademark of European Energy Exchange AG.

Flexible is better.® is a registered trademark of Axioma, Inc.